# A *REALLY* SIMPLE APPROXIMATION OF SMALLEST GRAMMAR

ARTUR JEŻ

ABSTRACT. In this paper we present a *really* simple linear-time algorithm constructing a context-free grammar of size $\mathcal{O}(g \log(N/g))$ for the input string, where $N$ is the size of the input string and $g$ the size of the optimal grammar generating this string. The algorithm works for arbitrary size alphabets, but the running time is linear assuming that the alphabet $\Sigma$ of the input string can be identified with numbers from $\{1, \ldots, N^c\}$ for some constant $c$. Algorithms with such an approximation guarantee and running time are known, however all of them were non-trivial and their analyses were involved. The here presented algorithm computes the LZ77 factorisation and transforms it in phases to a grammar. In each phase it maintains an LZ77-like factorisation of the word with at most $\ell$ factors as well as additional $\mathcal{O}(\ell)$ letters, where $\ell$ was the size of the original LZ77 factorisation. In one phase in a greedy way (by a left-to-right sweep and a help of the factorisation) we choose a set of pairs of consecutive letters to be replaced with new symbols, i.e. nonterminals of the constructed grammar. We choose at least 2/3 of the letters in the word and there are $\mathcal{O}(\ell)$ many different pairs among them. Hence there are $\mathcal{O}(\log N)$ phases, each of them introduces $\mathcal{O}(\ell)$ nonterminals to a grammar. A more precise analysis yields a bound $\mathcal{O}(\ell \log(N/\ell))$. As $\ell \leq g$, this yields the desired bound $\mathcal{O}(g \log(N/g))$.

## 1. INTRODUCTION

**Grammar based compression.** In the grammar-based compression text is represented by a context-free grammar (CFG) generating exactly one string. Such an approach was first considered by Rubin [20], though he did not mention CFGs explicitly. In general, the idea behind this approach is that a CFG can compactly represent the structure of the text, even if this structure is not apparent. Furthermore, the natural hierarchical definition of the context-free grammars makes such a representation suitable for algorithms, in which case the string operations can be performed on the compressed representation, without the need of the explicit decompression [5, 8, 13, 19, 6, 2].

While grammar-based compression was introduced with practical purposes in mind and the paradigm was used in several implementations [15, 14, 17], it also turned out to be very useful in more theoretical considerations. Intuitively, in many cases large data have relatively simple inductive definition, which results in a grammar representation of small size. On the other hand, it was already mentioned that the hierarchical structure of the CFGs allows operations directly on the compressed representation. A recent survey by Lohrey[16] gives a comprehensive description of several areas of theoretical computer science in which grammar-based compression was successfully applied.

The main drawback of the grammar-based compression is that producing the smallest CFG for a text is *intractable*: given a string $w$ and number $k$ it is NP-hard to decide whether there exist a CFG of size $k$ that generates $w$ [23]. Furthermore, the size of the smallest grammar for the input string cannot be approximated within some small constant factor [2].

**Previous approximation algorithms.** The first two algorithms with an approximation ratio $\mathcal{O}(\log(N/g))$ were developed simultaneously by Rytter [21] and Charikar et al. [2]. They followed a similar approach, we first present Rytter's approach as it is a bit easier to explain.

Rytter's algorithm [21] applies the LZ77 compression to the input string and then transforms the obtained LZ77 representation to an $\mathcal{O}(\ell \log(N/\ell))$ size grammar, where $\ell$ is the size of the LZ77 representation. It is easy to show that $\ell \leq g$ and as $f(x) = x \log(N/x)$ is increasing,

the bound $\mathcal{O}(g \log(N/g))$ on the size of the grammar follows (and so a bound $\mathcal{O}(\log(N/g))$ on approximation ratio). The crucial part of the construction is the requirement that the derivation tree of the intermediate constructed grammar satisfies the AVL condition. While enforcing this requirement is in fact easier than in the case of the AVL search trees (as the internal nodes do not store any data), it remains involved and non-trivial. Note that the final grammar for the input text is also AVL-balanced, which makes is suitable for later processing.

Charikar et al. [2] followed more or less the same path, with a different condition imposed on the grammar: it is required that its derivation tree is length-balanced, i.e. for a rule $X \to YZ$ the lengths of words generated by $Y$ and $Z$ are within a certain multiplicative constant factor from each other. For such trees efficient implementation of merging, splitting etc. operations were given (i.e. constructed from scratch) by the authors and so the same running time as in the case of the AVL grammars was obtained. Since all the operations are defined from scratch, the obtained algorithm is also quite involved and the analysis is even more non-trivial.

Sakamoto [22] proposed a different algorithm, based on RePair [15], which is one of the practically implemented and used algorithms for grammar-based compression. His algorithm iteratively replaces pairs of different letters and maximal repetitions of letters ($a^\ell$ is a *maximal repetition* if that cannot be extended by $a$ to either side). A special pairing of the letters was devised, so that it is 'synchronising': if $u$ has 2 disjoint occurrences in $w$, then those two occurrences can be represented as $u_1 u' u_2$, where $u_1, u_2 = \mathcal{O}(1)$, such that both occurrences of $u'$ in $w$ are paired and compressed in the same way. The analysis was based on considering the LZ77 representation of the text and proving that due to 'synchronisation' the factors of LZ77 are compressed very similarly as the text to which they refer. Constructing such a pairing is involved and the analysis non-trivial.

Recently, the author proposed another algorithm [9]. Similarly to the Sakamoto's algorithm it iteratively applied two local replacement rules (replacing pairs of different letters and replacing maximal repetitions of letters). Though the choice of pairs to be replaced was simpler, still the construction was involved. The main feature of the algorithm was its analysis based on the recompression technique, which allowed avoiding the connection of SLPs and LZ77 compression. This made it possible to generalise this approach also to grammars generating trees [10]. On the downside, the analysis is quite complex.

**Contribution of this paper.** We present a very simple algorithm together with a straightforward and natural analysis. It chooses the pairs to be replaced in the word during a left-to-right sweep and additionally using the information given by a LZ77 factorisation. We require that any pair that is chosen to be replaced is either inside a factor of length at least 2 or consists of two factors of length 1 and that the factor is paired in the same way as its definition. To this end we modify the LZ77 factorisation during the sweep. After the choice, the pairs are replaced and the new word inherits the factorisation from the original word. This procedure is repeated until a trivial word is obtained. To see that this is indeed a grammar construction, when the pair $ab$ is replaced by $c$ we create a rule $c \to ab$.

*Note on computational model.* The presented algorithm runs in linear time, assuming that we can compute the LZ77 factorisation in linear time. This can be done assuming that the letters of the input words can be sorted in linear time, which follows from a standard assumption that $\Sigma$ can be identified with a continues subset of natural numbers of size $\mathcal{O}(N^c)$ for some constant $c$ and the RadixSort can be performed on it. Note that such an assumption is needed for all currently known linear-time algorithms that attain the $\mathcal{O}(\log(N/g))$ approximation guarantee.

## 2. Notions

By $N$ we denote the size of the input word.

**LZ77 factorisation.** An LZ77 factorisation (called simply factorisation in the rest of the paper) of a word $w$ is a representation $w = f_1 f_2 \cdots f_\ell$, where each $f_i$ is either a single letter (called *free letter* in the following) or $f_i = w[j \mathinner{.\,.} j + |f_i| - 1]$ for some $j \le |f_1 \cdots f_{i-1}|$, in such a case $f$ is called a *factor* and $w[j \mathinner{.\,.} j + |f_i| - 1]$ is called the *definition* of this factor. We do not assume

that a factor has more than one letter though when we find such a factor we demote it to a free letter. The *size* of the LZ77 factorisation $f_1 f_2 \cdots f_\ell$ is $\ell$. There are several simple and efficient linear-time algorithms for computing the smallest LZ77 factorisation of a word [1, 3, 4, 7, 11, 18] and all of them rely on linear-time algorithm for computing the suffix array [12].

**SLP.** *Straight Line Programme* (*SLP*) is a CFG in the Chomsky normal form that generates a unique string. Without loss of generality we assume that nonterminals of an SLP are $X_1, \ldots, X_g$, each rule is either of the form $X_i \to a$ or $X_i \to X_j X_k$, where $j, k < i$. The *size* of the SLP is the number of its nonterminals (here: $g$).

The problem of finding smallest SLP generating the input word $w$ is NP-hard [23] and the size of the smallest grammar for the input word cannot be approximated within some small constant factor [2]. On the other hand, several algorithms with approximation ratio $\mathcal{O}(1 + \log(N/g))$, where $g$ is the size of the smallest grammar generating $w$, are known [2, 21, 22, 9]. Most of those constructions use the inequality $\ell \leq g$, where $\ell$ ($g$) is the size of the smallest LZ77 factorisation (grammar, respectively) for $w$ [21].

## 3. INTUITION

**Pairing.** Relaxing the Chomsky normal form, let us identify each nonterminal generating a single letter with this letter. Suppose that we already have an SLP for $w$. Consider the derivation tree for $w$ and the nodes that have only leaves as children (they correspond to nonterminals that have only letters on the right-hand side). Such nodes define a *pairing* on $w$, in which each letter is paired with one of the neighbouring letters (such pairing is of course a symmetric relation). Construction of the grammar can be naturally identified with iterative pairing: for a word $w_i$ we find a pairing, replace pairs of letters with 'fresh' letters (different occurrences of a pair $ab$ can be replaced with the same letter though this is not essential), obtaining $w_{i+1}$ and continue the process until a word $w_{i'}$ has only one letter. The fresh letters from all pairings are the nonterminals of the constructed SLP and its size is twice the number of different introduced letters. Our algorithm will find one such pairing using the LZ77 factorisation of a word.

**Creating a pairing.** Suppose that we are given a word $w$ and know its factorisation. We try the naive pairing: the first letter is paired with second, third with fourth and so on, see Fig. 1. If we now replace all pairs with new letters, we get a word that is 2 times shorter so $\log N$ such iterations give an SLP for $w$. However, in the worst case there are $|w|/2$ different pairs already in the first pairing and so we cannot give any better bound on the grammar size than $\mathcal{O}(N)$.

A better estimation uses the LZ77 factorisation. Let $w = f_1 f_2 \cdots f_\ell$ and consider a factor $f_i$. It is equal to $w[j \mathinner{.\,.} j + |f_i| - 1]$ and so all pairs occurring in $f_i$ already occur in $w[j \mathinner{.\,.} j + |f_i| - 1]$ unless the parity is different, i.e. $j$ and $|f_1 \cdots f_{i-1}| + 1$ are of different parity, see Fig. 1. We want to fix this: it seems a bad idea to change the pairing in $w[j \mathinner{.\,.} j + |f_i| - 1]$, so we change it in $f_i$: it is enough to shift the pairing by one letter, i.e. leave the first letter of $f_i$ unpaired and pair the rest as in $w[j + 1 \mathinner{.\,.} j + |f_i| - 1]$. Note that the last letter in the factor definition may be paired with the letter to the right, which may be impossible inside $f_i$, see Fig. 1. As a last observation note that since we alter each $f_i$, instead of creating a pairing at the beginning and modifying it we can create the pairing while scanning the word from left to right.

There is one issue: after the pairing we want to replace the pairs with fresh letters. This may make some of the the factor definitions improper: when $f_i$ is defined as $w[j \mathinner{.\,.} j + |f_i| - 1]$ it might be that $w[j]$ is paired with letter to the left. To avoid this situation, we replace the factor $f_i$ with $w[j] f_i'$ and change its definition to $w[j + 1 \mathinner{.\,.} j + |f_i| - 1]$. Similar operation may be needed at the end of the factor, see Fig. 1. This increases the size of the LZ77 factorisation, but the number of factors stays the same (i.e. only the number of free letters increases). Additionally, we pair two neighbouring free letters, whenever this is possible.

**Using a pairing.** When the appropriate pairing is created, we replace each pair with a new letter. If the pair is within a factor, we replace it with the same symbol as the corresponding pair in the definition of the factor. In this way only pairs that are formed from free letters may
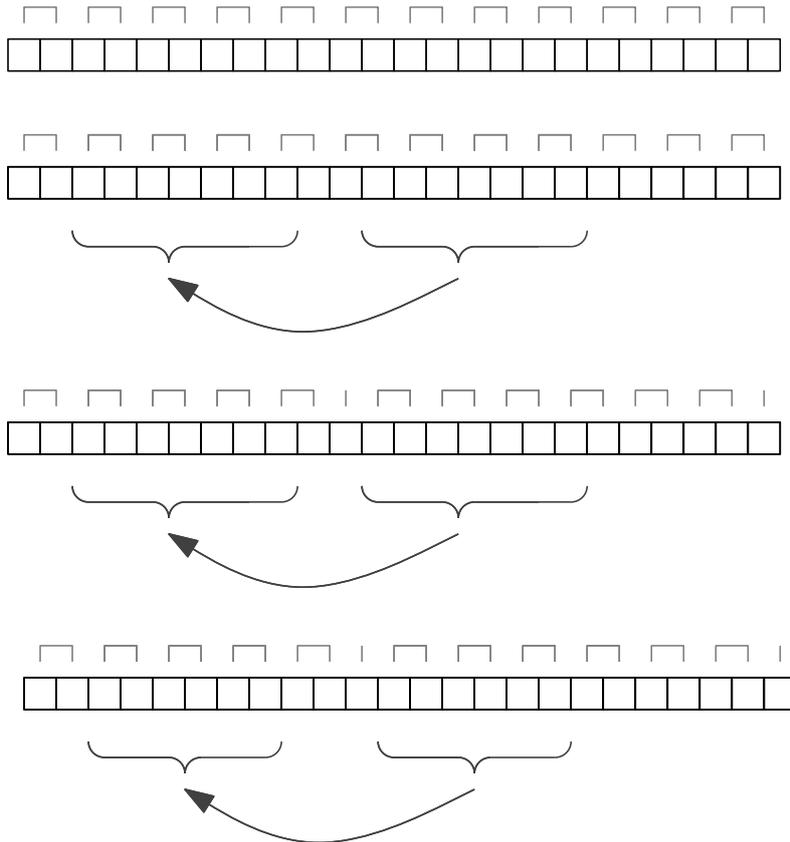
FIGURE 1. The pairings are presented over the words, the LZ77 factors are depicted below the words. On the top, the naive pairing is presented. On the second picture, the pairing is compared with the LZ77 factor as well as its definition; the factor and its definition are paired differently. On the third picture, we move the pairing so that it is consistent on the factor and its definition. This creates unpaired letters. On the bottom, we shorten the factor on the right, so that it ends with a whole pair.

contribute a fresh letter. As a result we obtain a new word together with a factorisation in which there are $\ell$ factors.

**Analysis.** The analysis is based on the observation that a factor $f_i$ is shortened by a constant fraction, so it takes part in $\log|f_i|$ phases and in each of them it introduces $\mathcal{O}(1)$ free letters. Hence the total number of free letters introduced to the word is $\mathcal{O}(\sum_{j=1}^{\ell} \log|f_i|) = \mathcal{O}(\ell \log(N/\ell))$ (which is shown in details later on). As creation of a rule decreases the number of free letters in the instance by at least 1, we obtain that this is also an upper bound on the size of the grammar.

## 4. The algorithm

**Stored data.** The word is represented as a table of letters. The table *start* stores the information about the beginnings of factors: $start[i] = j$ means that $i$ is the first letter of a factor and $j$ is the first letter of its definition; otherwise $start[i] = false$. Similarly *end* stores the information about the ends of factors: $end[i]$ is a bit flag, i.e. has value *true/false*, that tells whether $w[i]$ is the last letter of a factor.

When we replace the pairs with new letters, we reuse the same tables, overwriting from left to the right. Additionally, a table *newpos* stores the corresponding positions: $newpos[i] = j$ means that

- the letter on position $i$ was unpaired and $j$ is the position of the corresponding letter in the new word *or*

- the letter on position $i$ was paired with a letter to the right and the corresponding letter in the new word is on position $j$.

Lastly, *newpos*$[i]$ is not defined when position $i$ was the second element in the pair.

It is easy to see that the algorithm can be converted to lists and pointers instead of tables. (Though the RadixSort used in the LZ77 construction needs tables).

**Technical assumption.** Our algorithm makes a technical assumptions: a factor $f_i$ (of length at least 2) starting at position $j$ cannot have $start[j] = j - 1$, i.e. its definition is at least two positions to the left. This is verified and repaired while sweeping through $w$: if $w[j \mathinner{\ldotp\ldotp} j+|f_i|-1] = w[j-1 \mathinner{\ldotp\ldotp} j+|f_i|-2]$ then $f_i = a^{|f_i|}$ for some letter $a$. In such a case we split $f_i$: we make $w[j]$ a free letter and set $start[j+1] = j - 1$ (note that the latter essentially requires that indeed $f_i$ is a repetition of one letter).

**Pairing.** We are going to devise a pairing with the following (a bit technical) properties:

(P1) there are no two consecutive letters that are both unpaired;

(P2) the first two letters (last two letters) of any factor $f$ are paired with each other;

(P3) if $f = w[i \mathinner{\ldotp\ldotp} i + |f| - 1]$ has a definition $w[start[i] \mathinner{\ldotp\ldotp} start[i] + |f| - 1]$ then letters in $f$ and in $w[start[i] \mathinner{\ldotp\ldotp} start[i] + |f| - 1]$ are paired in the same way.

The pairing is found incrementally by a left-to-right scan through $w$: we read $w$ and when we are at letter $i$ we make sure that the word $w[1 \mathinner{\ldotp\ldotp} i]$ satisfies (P1)–(P3). To this end we not only devise the pairing but also modify the factorisation a bit (by replacing a factor $f$ with $af$ or by $fb$, where $a$ is the first and $b$ the last letter of $f$). If during the sweep some $f$ is shortened so that $|f| = 1$ then we demote it to a free letter.

The pairing is recorded in table: *pair*$[i]$ can be set to *first*, *second* or *none*, meaning that $w[i]$ is the first, second in the pair or it is unpaired, respectively.

**Creation of pairing.** We read $w$ from left to right, suppose that we are at position $i$. If we read $w[i]$ that is a free letter then we check, whether the previous letter is not paired. If so, then we pair them. Otherwise we continue to the next position.

If $i$ is a first letter of a factor, we first check whether the length of this factor is one; if so, we change $w[i]$ into a free letter. If the factor has definition only one position to the left (i.e. at $i-1$) then we split the factor: we make $w[i]$ a free letter and set $w[i+1]$ as a first letter of a factor with a definition starting at $i-1$. Otherwise we check whether $w[start[i]]$ is indeed the first letter of a pair. If not (i.e. it is a second letter of a pair or an unpaired letter) then we split the factor: we make $w[i]$ a free letter and $w[i+1]$ the beginning of a factor with a definition beginning at $start[i] + 1$ (note that this factor may have length 1); we view the factor beginning at $w[i+1]$ as a modified factor that used to begin at $w[i]$. If for any reason we turned $w[i]$ into a free letter, we re-read this letter, treating it accordingly. If $w[start[i]]$ is a first letter of a pair, we copy the pairing from the whole factor's definition to the factor starting at $i$.

When this is done we need to make sure that the factor indeed ends with a pair, i.e. that (P2) holds: if the last letter of a factor, say $w[i']$, is not the second in the pair. To this end we split the factor: we make $w[i']$ a free letter, clear $i'$'s pairing, decrease $i'$ by 1 and set the flag for $w[i']$ (making it the end of the new factor). We iterate it until the $w[i']$ is indeed a second letter of a factor. This is all formalised in Algorithm 1.

*Using the pairing.* When the pairing is done, we read the word $w$ again (from left to right) and replace the pairs by letters. We keep two indices: $i$, which is the pointer in the current word (pointing at the first unread letter) and $i'$, which is a pointer in the new word, pointing at the first free position. Additionally, when reading $i$ we store (in *newpos*$[i]$) the position of the corresponding letter in the new word, which is always $i'$.

If $w[i]$ is a first letter in a pair and this pair consists of two free letters, in the new word we add a fresh letter and move two letters to the right in $w$ (as well as one position in the new word). If $w[i]$ is unpaired and a free letter then we simply copy this letter to the new word, increasing both $i$ and $i'$ by 1. If $w[i]$ is first letter of a factor (and so also a first letter of a pair by (P2)), we copy the corresponding fragment of the new word (the first position is given by

---

**Algorithm 1** Pairing

---

1: $pair[1] \leftarrow none$
2: $i \leftarrow 2$
3: **while** $i \leq |w|$ **do**
4:     **if** $start[i]$ **then**                                    ▷ $w[i]$ is the first element of a factor
5:         **if** $end[i]$ **then**                                          ▷ This is one-letter factor
6:             $start[i] \leftarrow end[i] \leftarrow false$                          ▷ Turn it into a free letter
7:         **else if** $start[i] = i - 1$ **then**  ▷ The factor is $a^k$, its definition begins one position to
   the left
8:             $start[i + 1] \leftarrow i - 1$                             ▷ Move the definition of the factor
9:             $start[i] \leftarrow false$                                   ▷ Make $w[i]$ a free letter
10:         **else if** $pair[start[i]] \neq first$ **then**        ▷ The pairing of the definition of factor is bad
11:             $start[i + 1] \leftarrow start[i] + 1$                        ▷ Shorten the factor definition
12:             $start[i] \leftarrow false$                                 ▷ Make $w[i]$ a free letter
13:         **else**                                                        ▷ Good factor
14:             $j \leftarrow start[i]$                                   ▷ Factor's definition begins at $j$
15:             **repeat**                                       ▷ Copy the pairing from the factor definition
16:                 $pair[i] \leftarrow pair[j]$
17:                 $i \leftarrow i + 1, j \leftarrow j + 1$
18:             **until** $end[i - 1]$
19:             **while** $pair[i - 1] \neq second$ **do**                 ▷ Looking for a new end of the factor
20:                 $i \leftarrow i - 1$
21:                 $end[i - 1] = true, end[i] = false$                          ▷ Shortening of the factor
22:                 $pair[i] = none$                                         ▷ Clear the pairing
23:     **if** **not** $start[i]$ **then**                                      ▷ $w[i]$ a free letter
24:         **if** $pair[i - 1] = none$ **then**                      ▷ If previous letter is not paired
25:             $pair[i - 1] \leftarrow first, pair[i] \leftarrow second$                          ▷ Pair them
26:         **else**
27:             $pair[i] \leftarrow none$                               ▷ Leave the letter unpaired
28:     $i \leftarrow i + 1$

---

$newpos[start[i]])$, moving $i$ and $i'$ in parallel: $i'$ is always incremented by 1, while $i$ is moved by 2 when it reads a first letter of a pair and by 1 when it reads a free letter. Also, we store the new beginning and end of the factor in the new word: for a factor beginning at $i$ and ending at $i'$ we set $start[newpos[i]] = newpos[start[i]]$ and $end[newpos[i' - 1]] = true$ (note that $i' - 1$ and $i'$ are paired). Details are given in Algorithm 2.

**Algorithm.** TtoG first computes the LZ77 factorisation and then iteratively applies Pairing and PairReplacement, until a one-word letter is obtained.

## 5. Analysis

We begin the analysis with showing that indeed Pairing produces the desired pairing.

**Lemma 1.** Pairing *runs in linear time. It creates a proper factorisation and returns a pairing that satisfies (P1)–(P3) (for this new factorisation). When the current factorisation for the input word for* Pairing *has $m$ factors then* Pairing *creates at most $6m$ new free letters and the returned pairing has at most $m$ factors.*

*Proof.* For the running time analysis note that a single letter can be considered at most twice: once as a part of a factor and once as a free letter.

We show the second claim of the lemma by induction: at all time the stored factorisation is proper, furthermore, when we processed $w[1 . . i - 1]$ (i.e. we are at position $i$, note that we can go back in which case position gets unprocessed) then we have a partial pairing on $w[1 . . i - 1]$,

---

**Algorithm 2** PairReplacement

---

1:  $i \leftarrow i' \leftarrow 1$                                                $\triangleright$ $i'$ is the position corresponding to $i$ in the new word
2:  **while** $i \leq |w|$ **do**
3:      **if** $start[i]$ **then**                                        $\triangleright$ $w[i]$ is the first element of a factor
4:          $start[i'] \leftarrow j' \leftarrow newpos[start[i]]$          $\triangleright$ Factor in new word begins at the position
    corresponding to the beginning of the current factor
5:          $start[i] \leftarrow false$                                      $\triangleright$ Clearing obsolete information
6:          **repeat**
7:              $newpos[i] \leftarrow i'$                                  $\triangleright$ Position corresponding to $i$
8:              $w[i'] \leftarrow w[j']$                        $\triangleright$ Copy the letter according to new factorisation
9:              $i' \leftarrow i' + 1$, $j' \leftarrow j' + 1$
10:             **if** $pair[i] = first$ **then**
11:                 $i \leftarrow i + 2$                                  $\triangleright$ We move left by the whole pair
12:             **else**
13:                 $i \leftarrow i + 1$                          $\triangleright$ We move left by the unpaired letter
14:          **until** $end[i - 1]$                              $\triangleright$ We processed the whole factor
15:          $end[i' - 1] \leftarrow true$                                $\triangleright$ End in the new word
16:          $end[i - 1] \leftarrow false$                          $\triangleright$ Clearing obsolete information
17:     **if** **not** $start[i]$ **then**                                        $\triangleright$ $w[i]$ a free letter
18:         $newpos[i] \leftarrow i'$
19:         **if** $pair[i] = none$ **then**
20:             $w[i'] \leftarrow w[i]$                              $\triangleright$ We copy the unpaired letter
21:             $i \leftarrow i + 1$, $i' \leftarrow i' + 1$                $\triangleright$ We move by this letter to the right
22:         **else**
23:             $w[i'] \leftarrow$ fresh letter                $\triangleright$ Paired free letters are replaced by a fresh letter
24:             $i \leftarrow i + 2$, $i' \leftarrow i' + 1$              $\triangleright$ We move to the right by the whole pair

---

**Algorithm 3** TtoG

---

1:  compute LZ77 factorisation of $w$
2:  **while** $|w| > 1$ **do**
3:      compute a pairing of $w$ using Pairing
4:      replace the pairs using PairReplacement
5:  output the constructed grammar

---

which differs from the pairing only in the fact that the position $i - 1$ may be assigned as first in the pair and $i$ is not yet paired. This partial pairing satisfies (P1)–(P3) restricted to $w[1 \mathinner{.\,.} i - 1]$.

*Factorisation.* We first show that after considering $i$ the modified factorisation is proper.

If in line 5 we have $start[i] = end[i]$ then $w[i]$ is a one-letter factor and so after replacing it with a free letter the factorisation stays proper. Observe now that the verification in line 5 ensures that in each other case considered in lines 7–22 we deal with factors of length at least 2.

The modifications of the factorisation in line 9 results in a proper factorisation: the change is applied only when $start[i] = i - 1$, in which case $w[i \mathinner{.\,.} i + |f| - 1] = w[i - 1 \mathinner{.\,.} i + |f| - 2]$, which implies that $f = a^{|f|}$, where $a = w[i]$. Since $|f| \geq 2$ (by case assumption), in such a case $w[i + 1 \mathinner{.\,.} i + |f| - 1] = w[i - 1 \mathinner{.\,.} i + |f| - 3]$ so we can split the factor $f = w[i \mathinner{.\,.} i + |f| - 1]$ to $w[i]$ and a factor $w[i + 1 \mathinner{.\,.} i + |f| - 1]$ which is defined as $w[i - 1 \mathinner{.\,.} i + |f| - 3]$ ($f$ had at least two letters, so after the modification it has at least 1 letter).

In line 11 we shorten the factor by one letter (and create a free letter), as the factor had at least two letters, so the factorisation remains proper.

Concerning the symmetric shortening in line 20, it leaves a proper factorisation (as in case of line 11), as long as we do not move $i$ before the beginning of the factor. However, observe that when we reach line 16 this means that the factor beginning at $i'$ has length at least 2,

$start[i'] < i' - 1$ and $pair[start[i']] = first$. Thus $start[i'] + 1 < i'$ and so by induction assumption we already made a pairing for it. Since $start[i']$ is assigned $first$, $start[i'] + 1$ is assigned $second$. So $i' + 1$ is assigned $second$ as well. Since the end of the factor is at position $i \geq i' + 1$, in our search for element marked with $second$ at positions $i$, $i - 1$, ... we cannot move to the left more than to $i' + 1$. Thus the factor remains (and has at least 2 letters).

*Pairing.* We show that indeed we have a partial pairing. Firstly, if $i$ is decreased, then as a result we get a partial pairing: the only nontrivial case is when $i - 1$ and $i$ were paired then $i - 1$ is assigned as the first element in the pair but it has no corresponding element, which is allowed in the partial pairing. If $i$ is increased then we need to make sure if $i - 1$ is assigned as a first element in a pair then $i$ will be assigned as the second one (or the pairing is cleared). Note that $i - 1$ can be assigned in this way only when it is part of the factor, i.e. it gets the same status as some $j$. If $i$ is also part of the same factor, then it is assigned the status of $j + 1$, which by inductive assumption is paired with $j$, so is the second element in the pair. In the remaining case, if $i - 1$ was the last element of the factor then in loop in line 19 we decrease $i$ and so unprocess $i - 1$ (in particular, we clear its pairing).

For (P2) observe that for the first two letters it is explicitly verified in line 10. Similarly, for the second part of (P2): we shorten the last factor in line 19 (ending at $i$) until $pair[i] = second$. We already shown that pairing is defined for $w[1 \ldots i]$ and when $i$ is assigned $second$ then $i - 1$ is assigned $first$, as claimed.

Condition (P3) is explicitly enforced in loop in line 15, in which we copy the pairing from the definition of the factor.

Suppose that (P1) does not hold for $i - 1, i$, i.e. they are both unpaired after processing $i$. It cannot be that they are both within the same factor, as then the corresponding $w[j - 1]$ and $w[j]$ in the definition of the factor are also unpaired, by (P3), which contradicts the induction assumption. Similarly, it cannot be that one of them is in a factor and the other outside this factor, as by (P2) (which holds for $w[1 \ldots i]$) a factor begins and ends with two paired letters. So they are both free letters. But then we needed to pass line 23 for $i$ and both $w[i-1]$ and $w[i]$ were free at that time, which means that they should have been paired at that point, contradiction.

To see the third claim of the lemma, i.e. the bound on the number of new free letters, fix a factor $f$ that begins at position $i$. When it is modified, we identify the obtained factor with $f$ (which in particular shows that the number of factors does not increase). We show that it creates at most 6 new free letters in this phase.

If at any point $start[i]$ and $end[i]$, i.e. the factor has only one letter, then it is replaced with a free letter and afterwards cannot introduce any free letters (as $f$ is no longer there). Hence at most one free letter is introduced by $f$ due to condition in line 5.

If $start[i] = i - 1$ then it creates one free letter inside condition in line 7. It cannot introduce another free letter in this way (in this phase), as afterwards $start[i + 1] = i - 1$ and there is no way to decrease this distance (in this phase).

We show that condition in line 10, i.e. that the first letter of the factor definition is not the first in the pair, holds at most twice for a fixed factor $f$ in a phase. Since we set $j = start[i]$ and increase both $i$ and $j$ by 1 until $pair[j] = first$, this can be viewed as searching for the smallest position $j' \geq j$ that is first in a pair and we claim that $j' \leq j + 2$. On the high-level, this should hold because (P1) holds for $w[1 \ldots i - 1]$, and so among three consecutive letters there is at least one that is the first element in the pair. However, the situation is a bit more complicated, as some pairing may change during the search.[1]

Consider first the main case, in which $j \leq i - 3$. Then the elements at position $j$, $j + 1$, $j + 2 \leq i - 1$ have already assigned pairings and so at least one of them is assigned as a first element of some pair. The only way to change the pairing from $first$ to some other is in loop in line 19. However, we can go to this loop only after condition in line 10 fails, which implies that it holds at most twice (i.e. for at most two other among $j$, $j + 1$, $j + 2$).

As the case in which $start[i] = i - 1$ is excluded by the case assumption of the algorithm, the remaining case is $j = i - 2$. As in the previous argument, we consider the letters whose pairing

---

[1]In particular, this could fail if we allowed that $start[i] = i - 1$, so some care is needed.

are known, i.e. $w[i-2]$ and $w[i-1]$. If any of them is a first letter in a pair, we are done (as in the previous case). As (P1) holds for them, the only remaining possibility is that $w[i-2]$ is a second letter in a pair and $w[i-1]$ is unpaired. Then when we consider $i$ in line 10 it is made free. When we consider $w[i]$ in line 23 (re-reading it as a free letter), it is paired with $i-1$. Hence when we read $i+1$ (the new first letter in the factor) its definition $(i-1)$ is a first letter in a pair, as claimed.

Similar analysis can be applied to the last letter of a factor. So, as claimed, one factor can introduce at most 6 free letters in a single phase.

Finally, it is left to show that when we processed the whole $w$ then we have a proper factorisation and a pairing satisfying (P1)–(P3). From the inductive proof it follows that the kept factorisation is proper and the partial pairing satisfies (P1)–(P3) for the whole word. So it is enough to show that this partial pairing is a pairing, i.e. that the last letter of $w$ is not assigned as a first element of a pair. Consider, whether it is in a factor or a free letter. If it is in a factor then clearly it is the last element of the factor and so by (P2) it is the second element in a pair. If it is a free letter observe that we only pair free letters in line 25, which means that it is paired with the letter on the next position, contradiction. □

Now, we show that when we have a pairing satisfying (P1)–(P3) (so in particular the one provided by Pairing is fine, but it can be any other pairing satisfying (P1)–(P3)) then PairReplacement creates a word $w'$ out of $w$ together with a factorisation.

**Lemma 2.** *When a pairing satisfies (P1)–(P3) then* PairReplacement *runs in linear time and returns a word $w'$ together with a factorisation; $|w'| \leq \frac{2|w|+1}{3}$ and the factorisation of $w'$ has the same number of factors as the factorisation of $w$. If $p$ fresh letters were introduced then $w'$ has $p$ less free letters than $w$.*

*Proof.* The running time is obvious as we make one scan through $w$.

Firstly, we show that when we erase the information about beginnings and ends of factors of $w$ we do not erase the newly created information for $w'$. To this end it is enough to show that $i > i'$ in such situation. Whenever $i'$ is incremented, $i$ is incremented by at least the same amount, so it is enough to show that $i > i'$ when $i$ is the first letter of the first factor, in other words, there is at least one pair before the first factor. By (P1) there is a pair within first three positions of the factor. If the pair is at positions 1, 2 then by (P2) the factor begins at position 3 or later and we are done. If the pair is at positions 2, 3 then by (P2) the factor begins at position 4 or later or at position 2; however, the latter case implies that the factor definition is 1 to the left of the factor, which is excluded by Pairing.

Concerning the size of the produced word, by (P1) each unpaired letter (perhaps except the last letter of $w$) is followed by a pair. Thus, at least $\frac{1}{3}(|w|-1)$ letters are removed from $w$, which yields the claim.

Concerning the factorisation of $w'$, observe that by an easy induction it can be shown that for each $i$ the $newpos[i]$ is

- undefined, when $i$ is second in a pair *or*
- is the position of the corresponding letter in $w'$.

Now, consider any factor $f$ in $w$ with a definition $w[j \ldots j + |f| - 1]$. By (P2) both the first and the last two letters of $f$ are paired and by (P3) pairing of $f$ is the same as the pairing of its definition. So it is enough to copy the letters in $w'$ corresponding to $w[j \ldots j + |f| - 1]$, i.e. beginning with $newpos[j]$, which is what the algorithm does. When we consider a free letter, if it is unpaired, it should be copied (as it is not replaced), and when it is paired, the pair can be replaced with a fresh letter; in both cases the corresponding letter in the new word should be free. And the algorithm does that.

Concerning the number of fresh letters introduced, suppose that $ab$ is replaced with $c$. If $ab$ is within some factor $f$ then we use for the replacement the same letter as we use in the factor definition and so no new fresh letter is introduced. If both this $a$ and $b$ are free letters then each such a pair contributes one fresh letter. And those two free letters are replaced with one free letter, hence the number of free letters decreases by 1. The last possibility is that one letter

from $ab$ comes from a factor and the other from outside this factor, but this contradicts (P2) that a factor begins and ends with a pair. $\qquad\square$

Using the two lemmata we can give the proof.

**Theorem 1.** TtoG *runs in linear time and returns an SLP of size* $\mathcal{O}(\ell + \ell \log(N/\ell))$. *Thus its approximation ration is* $\mathcal{O}(1 + \log(N/g))$, *where* $g$ *is the size of the optimal grammar.*

*Proof.* Due to Lemma 2 each introduction of a fresh letter reduces the number of free letters by 1. Thus to bound the number of different introduced letters it is enough to estimate the number of created free letters. In the initial LZ77 factorisation there are at most $\ell$ of them. For the free letters created during the Pairing let us fix a factor $f$ of the original factorisation and estimate how many free letters it created. Due to (P1) the length of $f$ drops by a constant fraction in each phase and so it will take part in $\mathcal{O}(\log|f|)$ phases. In each phase it can introduce at most 6 free letters, by Lemma 1. So $\mathcal{O}(\sum_{i=1}^{\ell} \log|f_i|)$ free letters were introduced to the word during all phases. Consider $\sum_{i=1}^{\ell} \log|f_i|$ under the constraint $\sum_{i=1}^{\ell} |f_i| \leq N$. As function $h(x) = \log(x)$ is concave, we conclude that this is maximised for all $f_i$ being equal to $N/\ell$. Hence the number of nonterminals in the grammar introduced in this way is $\mathcal{O}(\ell \log(N/\ell))$. Adding the $\ell$ for the free letters in the LZ77 factorisations yields the claim.

Concerning the running time, the creation of the LZ77 factorisation takes linear time [1, 3, 4, 7, 11, 18]. In each phase the pairing and replacement of pairs takes linear time in the length of the current word. Thanks to (P1) the length of such a word is reduced by a constant fraction in each phase, hence the total running time is linear. $\qquad\square$

## References

[1] Anisa Al-Hafeedh, Maxime Crochemore, Lucian Ilie, Evguenia Kopylova, William F. Smyth, German Tischler, and Munina Yusufu. A comparison of index-based Lempel-Ziv LZ77 factorization algorithms. *ACM Comput. Surv.*, 45(1):5, 2012.

[2] Moses Charikar, Eric Lehman, Ding Liu, Rina Panigrahy, Manoj Prabhakaran, Amit Sahai, and Abhi Shelat. The smallest grammar problem. *IEEE Transactions on Information Theory*, 51(7):2554–2576, 2005.

[3] Gang Chen, Simon J. Puglisi, and William F. Smyth. Fast and practical algorithms for computing all the runs in a string. In Bin Ma and Kaizhong Zhang, editors, *CPM*, volume 4580 of *LNCS*, pages 307–315. Springer, 2007.

[4] Maxime Crochemore, Lucian Ilie, and William F. Smyth. A simple algorithm for computing the Lempel Ziv factorization. In *DCC*, pages 482–488. IEEE Computer Society, 2008.

[5] Paweł Gawrychowski. Pattern matching in Lempel-Ziv compressed strings: fast, simple, and deterministic. In Camil Demetrescu and Magnús M. Halldórsson, editors, *ESA*, volume 6942 of *LNCS*, pages 421–432. Springer, 2011.

[6] Leszek Gąsieniec, Marek Karpiński, Wojciech Plandowski, and Wojciech Rytter. Efficient algorithms for Lempel-Ziv encoding. In Rolf G. Karlsson and Andrzej Lingas, editors, *SWAT*, volume 1097 of *LNCS*, pages 392–403. Springer, 1996.

[7] Keisuke Goto and Hideo Bannai. Simpler and faster Lempel Ziv factorization. In Ali Bilgin, Michael W. Marcellin, Joan Serra-Sagristà, and James A. Storer, editors, *DCC*, pages 133–142. IEEE, 2013.

[8] Artur Jeż. Faster fully compressed pattern matching by recompression. In Artur Czumaj, Kurt Mehlhorn, Andrew Pitts, and Roger Wattenhofer, editors, *ICALP (1)*, volume 7391 of *LNCS*, pages 533–544. Springer, 2012.

[9] Artur Jeż. Approximation of grammar-based compression via recompression. In Johannes Fischer and Peter Sanders, editors, *CPM*, volume 7922 of *LNCS*, pages 165–176. Springer, 2013. full version at http://arxiv.org/abs/1301.5842.

[10] Artur Jeż and Markus Lohrey. Approximation of smallest linear tree grammar. In Ernst W. Mayr and Natacha Portier, editors, *STACS*, volume 25 of *LIPIcs*, pages 445–457. Schloss Dagstuhl — Leibniz-Zentrum fuer Informatik, 2014.

[11] Juha Kärkkäinen, Dominik Kempa, and Simon J. Puglisi. Linear time Lempel-Ziv factorization: Simple, fast, small. In Johannes Fischer and Peter Sanders, editors, *CPM*, volume 7922 of *LNCS*, pages 189–200. Springer, 2013.

[12] Juha Kärkkäinen, Peter Sanders, and Stefan Burkhardt. Linear work suffix array construction. *J. ACM*, 53(6):918–936, 2006.

[13] Marek Karpiński, Wojciech Rytter, and Ayumi Shinohara. Pattern-matching for strings with short descriptions. In *CPM*, pages 205–214, 1995.

[14] John C. Kieffer and En-Hui Yang. Sequential codes, lossless compression of individual sequences, and kolmogorov complexity. *IEEE Transactions on Information Theory*, 42(1):29–39, 1996.

[15] N. Jesper Larsson and Alistair Moffat. Offline dictionary-based compression. In *Data Compression Conference*, pages 296–305. IEEE Computer Society, 1999.

[16] Markus Lohrey. Algorithmics on SLP-compressed strings: A survey. *Groups Complexity Cryptology*, 4(2):241–299, 2012.

[17] Craig G. Nevill-Manning and Ian H. Witten. Identifying hierarchical strcture in sequences: A linear-time algorithm. *J. Artif. Intell. Res. (JAIR)*, 7:67–82, 1997.

[18] Enno Ohlebusch and Simon Gog. Lempel-Ziv factorization revisited. In Raffaele Giancarlo and Giovanni Manzini, editors, *CPM*, volume 6661 of *LNCS*, pages 15–26. Springer, 2011.

[19] Wojciech Plandowski. Testing equivalence of morphisms on context-free languages. In Jan van Leeuwen, editor, *ESA*, volume 855 of *LNCS*, pages 460–470. Springer, 1994.

[20] Frank Rubin. Experiments in text file compression. *Commun. ACM*, 19(11):617–623, 1976.

[21] Wojciech Rytter. Application of Lempel-Ziv factorization to the approximation of grammar-based compression. *Theor. Comput. Sci.*, 302(1-3):211–222, 2003.

[22] Hiroshi Sakamoto. A fully linear-time approximation algorithm for grammar-based compression. *J. Discrete Algorithms*, 3(2-4):416–430, 2005.

[23] James A. Storer and Thomas G. Szymanski. The macro model for data compression. In Richard J. Lipton, Walter A. Burkhard, Walter J. Savitch, Emily P. Friedman, and Alfred V. Aho, editors, *STOC*, pages 30–39. ACM, 1978.

Max Planck Institute für Informatik,, Campus E1 4, DE-66123 Saarbrücken, Germany, and Institute of Computer Science, University of Wrocław, ul. Joliot-Curie 15, 50-383 Wrocław, Poland,
AJE@CS.UNI.WROC.PL